

## 英語ニュースの機械翻訳

浦谷 則好 相沢 輝昭

NHK放送技術研究所

NHKでは衛星放送で英語ニュースのテロップ作成に英日機械翻訳システムを試用している。ニュース文は扱う分野が広い、複雑な固有名詞表現や数量表現が多い等の特徴がある。実際に機械翻訳で使用された英語ニュース文を分析したところ、1文当たり11.0語で、現在時制の平叙文が多いこと、無生物主語の文が多いこと等が判明した。また、機械翻訳での解析成功率は64.5%であった。衛星放送のニュース文を対象としていることによって生ずる点を中心に我々のシステムの問題点についても報告する。現在、英語ニュースデータベースを構築し、これを用いて辞書・文法の見直しを進め翻訳精度の向上を図っている。

## English-Japanese Machine Translation for News

Noriyoshi Uratani, Teruaki Aizawa

NHK Science and Research laboratories

1-10-11 Kinuta, Setagaya-ku, Tokyo, 157 Japan

NHK has been trying a machine translation system from English to Japanese for producing subtitles for news. News sentences have complex proper nouns and numeral expressions. We analyzed the English news sentences which were translated by our system. Most of them were declarative sentences with the present tense whose average length was 11.0 words. The success rate of translation was 64.5 percent. The problems of our system are described around the use in TV news. A database for AP news is being constructed. We are going to make use of it to improve the grammar and the dictionary of our MT system.

### 1. はじめに

機械翻訳システムは多くの電算機メーカー、ソフトウェアハウス、各種研究機関などにおいて研究開発が進められており、既に販売されているものも少なくはない。しかし、その多くは製品マニュアル、技術文献を対象にしており、翻訳精度も十分なものではない。

NHKでは、ニュースという広い分野を対象として、放送における機械翻訳技術の実用化を目指して、昭和63年度から英日機械翻訳の研究を行っている。現在、すでに衛星放送でテロップ（字幕）作成に試験的に使用しつつ<sup>1)</sup>、問題点の抽出、その改善を図っている。

以下、衛星放送における機械翻訳の試用状況や問題点を具体的に例示した後、今後の改良の進め方を述べる。また、翻訳精度向上のために構築している英語ニュースデータベースにもふれる。

### 2. 衛星放送における機械翻訳システムの試用状況

NHKにおける翻訳の需要は多い。例えば、衛星第1放送の「ワールドニュース」では、英・仏・独・伊・露・韓・中の各言語で入って来る外国語ニュースに迅速に日本語テロップを付けて放送しなければならない。この翻訳作業は時差の関係で深夜作業となることも多く、運用上の大きなネックとなっている。

機械翻訳システムは、この部分の将来における省力化を目的に導入され、昨年8月から英語ニュースの日本語テロップの作成に試用されている。ニュースの聴取から英文の要約までは人手で行う。テロップは通常15字程度であり、最大でも30字なので、この段階で相当大幅な要約が必要となる。この際、併せて英文の単純化などの前編集もなされる。この要約英文の日本語への翻訳を機械翻訳システムが担当する。翻訳された結果には後編集がなされ、日本語テロップが作成され、画面に合わせて送られる。(写真1)

放送用英日機械翻訳システムの現在の諸元は表1のようになっている。翻訳の基本ソフトウェアとしては既存のもの<sup>2)</sup>を導入した。

### 3. 英語ニュースの機械翻訳の分析結果

衛星第1放送「ワールドニュース」のために実際に機械翻訳に供された英語ニュース文(45日分)の分



写真1 実際の放送画面

表1 システム諸元

翻訳方式	トランスファ方式
辞書の語数	基本語4万語 専門語5.5万語
構文解析規則	約3000
ハードウェア	SUN3/260
使用言語	C
翻訳速度	約2万語/時

析を行った。結果を以下に示す。

#### 3.1 「ワールドニュース」文の特徴

対象ニュース文の各種統計は以下ようになった。

総文数	1,393
総語数	15,422 (ピリオド, コンマ除く)
・語数/文	11.0
・動詞/文	1.9
異なり語数	3,389 (数値を除く)
・固有名詞	291
・人名	120

文型は平叙文が1,088(78%)で最も多く、次はタイトル等の無動詞文(157:11%)であった。時制に関しては、主文が現在のものが902であり、補文でも現在時制が271と過去、未来に比して圧倒的に多く使用されていた。これは、画面と文が連

動しているテロップ文の特徴であろう。また、受動態の使用は134であった。本動詞以外の動詞の主な使われ方は、不定詞257、進行形112、完了形109であった。

その他では、無生物を主語とする文が501と多いことが目立った。これは翻訳対象としたニュースが政治経済が少なく、市民社会のトピック的なものが多かったためだろうと考えられる。また、従属節は221、AND並列は220、関係節は181出現し、これも比較的多いと言えよう。

### 3.2 システムの翻訳レベルの現状

「ワールドニュース」の1,393文に対する機械翻訳結果は次の通りである。

全文数	1,393
解析に成功	898 (64.5%)
第1候補	698 (77.7%)
第2候補	121
第3候補	70
解析に失敗	495
スベルミス	139 (28.1%)
非文	12

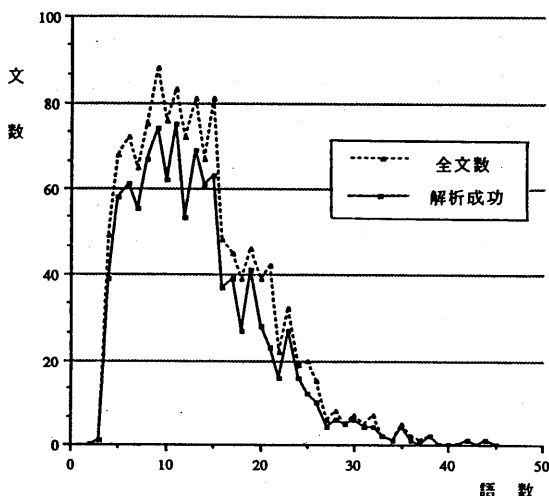


図1 語数と解析成功数との関係

文の語数と解析成功文の数との関係を図1に示す。

この結果から多くの候補文から最適の文を選択する機構(ウェイト機構)はおおよそうまく機能していると言えよう。解析に失敗したものの約30%は原文のスベルや非文が原因であり、それを除けば解析成功率は72.3%(898/1242)となる。

代表的な翻訳例を次に挙げる。

This is a new program exploring solutions to the most pressing environmental and social problems of our time.

㊦これは我々の時間の最も緊急の環境の、そして社会的問題の解を探索している新しいプログラムである。  
㊧並列句のスコープも正しく認定している。

Nearly 600 British soldiers and 300 police have been killed.

㊨ほとんど600イギリスの軍人と300警察は死んだ。

The Cheyenne Indians claim their ancestors fled to west before Columbus landed because prophecy warned them to flee the coming of a white destruction.

㊩予言がそれらに白い破壊の接近を避けなさいと警告したので、Columbusが上陸する前に、それらの先祖が西へ逃げたとCheyenneインディアンは主張する。

He plays easy listenig music as they munch on celery stocks.

㊪それらがセロリ株でむしゃむしゃ食べるように、彼はイージーリスニング音楽を演奏する。

㊫easy listenig musicが辞書に登録されているのでうまく訳されている。これが辞書に無いと「容易な聞いている音楽」と訳されて何のことか分からなくなる。

次に、我々の機械翻訳システムの問題点を、いくつかの翻訳結果をもとに、以下に例示する。ただし、品詞認定、訳語選択、to不定詞の訳し分け等どの翻

訳システムにも共通する問題は除外した。

★衛星放送ニュースを対象としたことから起こった問題点

◆画面を見ないと文意が分からないものがある。

They slither like a slity toves

⑩彼らは slity toves を slither 好む。

⑪この文はニュースキャスターがslitherからの連想で「鏡の国のアリス」のJabberwockyの詩の一節を思わず口ずさんだもの。この文だけだとミススペルとしか思えない。文脈処理を考えても対処不能と思われる。なお、slitherは未登録語であり、品詞認定を誤っている。

◆対象とする分野が広く、かつ個別には専門的。

(しかも学問的な専門分野と異なりトピック的)

Free-lance writer, Chuck Young, is trying to do purple haze just like Jimi Hendrix did it.

⑩フリーの作家 (Chuck Young) はJimi Hendrixのようなもやだけがそれにした紫色を行うことを…

⑪purple haze が曲名(「紫のけむり」)であると知らないで正しく訳せない。

◆間をカンマで表わすと、通常の構文解析が失敗。

"The new Tears For Fears album is called, "the seeds of love."

⑩新しいTears For Fears アルバムが呼ばれる ||, || 「愛の種」

⑪カンマが無ければもちろん正訳が得られる。

★前編集が不十分なために生じた問題点

◆ミススペルに弱い。

They began to disappear as land values shut up.

⑩彼らは土地価値が締まるように見えなくなり始めた。

⑪shut up は shot up (急騰する) の誤り。

スペルチェック機能があってもユーザは使ってくれと限らない!

◆口語的な表現が多い。

"This organized crime element, Russian emigrants, they kill anybody."

⑩「この組織された罪元素 || (ロシアの移民) 彼らは誰も殺す。」 ||

◆省略に弱い。

"Why not buy a motorized wheel chair?"

⑩「Whyが自動車化された車輪いすを買わない || ?」

◆文の分割に対する制約が厳しい。

While Irish fighters and Irish civilians, Catholic or Protestant, who have died, were 2000.

⑩アイルランドのが死んだ武人かアイルランドの一般国民かカトリックかプロテスタントは2000であった。

⑪従属節だけの文となっている。主文と一緒にすると正しい訳が得られた。

★その他

◆複合語、イディオムの登録が不十分

light house 誤訳：光家 正訳：灯台

air line 誤訳：空気ライン 正訳：航空会社

be scheduled to 誤訳：表に組み入れられている

正訳：することになっている

◆複雑な数量表現への対処が不十分

A family of four can see a show for under 10 ...

⑩4の1人の家族が、ショーを見ることができ…

この他にも係り先の認定、並列句のスコープ解釈、同格挿入句の解釈、分節構造の認定等の誤る例が見られる。しかし、これらは頻度も多くなく、また他の機械翻訳システムに比して劣っているとも考えられないので割愛した。さらに、助詞等の不適切な使用によって適切な日本語訳とならない例も多いが、今回は構文解析レベルを中心に報告したいのでこれも省く。

4. 翻訳精度改善のための英語ニュースデータベース

我々のシステムは放送用の機械翻訳システムとして、まだ十分とは言えない。ユーザーインターフェイスの改良としては、引き続き翻訳操作の改善や新語登録機能の充実などを図っていくつもりである。また機械翻訳システムと放送システムとのリンクもこれからの課

題である。テロップ送出装置とのオンライン接続など様々な形で機械翻訳を活用すべく検討を進めている。

3章で見てきたように我々のシステムの翻訳性能もまだ十分とは言えない。そこで、ニュース文では特徴的な固有名詞表現や数量表現の解析強化を先行して進める<sup>3)</sup>とともに、以下のような改善を進めている。

#### 4. 1 辞書の充実

機械翻訳システムは2種類の辞書を使い分けるようになっている。1つは、各語について詳細な記述を持つ基本語辞書であり、もう1つは基本的に英日対訳のみからなる専門語辞書である。

基本語辞書は当初2万語であったが、学習研究社の「学研英和電子辞書」を追加して4万語に拡張している。1つの見出し語には複数の訳語が対応しており、訳語の選択には文型パターンや意味マーカ―などを用いている。しかし、現状では特殊な語を除けば意味マーカ―は付与されていない。そこで、約80項目からなる意味マーカ―体系を構成し、各語への付与を進めている。<sup>4)</sup>

一方、専門語の方は、試用開始当初はゼロの状態であった。そこで、次の3種類の冊子型辞書を用いて充実に回り、現在5.5万語になっている。<sup>5)</sup>

三省堂『ニューズ英語辞典』

小学館『最新英語情報辞典』

集英社『情報・知識 i m i d a s 』

#### 4. 2 英語ニュースデータベースの構築

機械翻訳システムの翻訳精度の向上のためには、対象とする実際のデータに即して、辞書や文法をチューニングしていく必要がある。このため、我々はA P電を中心にニュースの実データの収集を行っている。A P電は昨年2月以降、A P通信社から直接電話線を紹介して入電しており、パソコンを用いて自動収集している。1日当たり約350記事(約700KB)が得られている。これを、効率良く利用するためにはデータベース化は不可欠である。

この英語ニュースデータベースの利用目的としては一応以下のものを想定している。

1) ニュース文そのものの検索

2) ニュース文からの専門語・慣用句(固有名詞表現、数値表現を含む)の抽出

3) ニュース文からの表現パターンの抽出

4) 辞書・文法の評価

データベース上のデータ表現としては、データの管理、保守、各種の処理を考えれば、関係データベースやオブジェクト表現が有利であろうが、データの分析(特にマイクロな分析)の手法が固まっていない現在ではむしろできるだけ生データに近い形にして後の処理に制約を加えたくないと考えた。そこで、単に記事毎にヘッダー、本文、トレイラを分離し、図2のように各記事のヘッダー、トレイラの部分と本文部分をまとめて格納することとした。図2で、“ $\$$ ”はヘッダー部とトレイラ部を分けるデリミタであり、“##”は記事本文やヘッダー・トレイラ部の各々を分けるデリミタである。さらに、データの先頭には記事毎のヘッダー・トレイラ部へのポインターと本文部へのポインターを付加した。A P電は日毎に、ワールドニュースは月毎にこのデータ形式に変換することにした。このような簡単な形式にしたので、原データに対してわずか1%の容量増で済ませることができた。自動収集されたA P電は変換コマンドによって自動的にEWS上のデータベースに取り込んでいるが、回線ノイズ等によってデータの欠落や、不正データの混入が起こることがあり、この場合はデータチェック等の支援コマンドを用いて人手で修正してデータベース化している。現在(90.6)A P電のデータは246MB格納されている。処理の基本となる文字列の検索は当所で開発したFAST法<sup>6)</sup>を組み込んで高速に実行している。

本格的なデータベースの利用に先立って、A P電データの予備調査を行った。A P電は1日に約350記事(約2KB/記事)送られてくるが、この中には要約記事、改作記事等の重複記事が多く含まれている。このままで各種統計をとっても有用なデータにはなれない。さらに、相場記事、スポーツ記事は一般ニュースとは別個に扱いたい。そこで、47個のキーワードを設定し、これをヘッダー中に含む記事をFAST法で探索し、不要と思われる記事を統計の対象から除外した。

1989年4月~90年3月(データの欠落がある

ため実質約290日分)の約1600万語について調査してみた。異なり語数は約17万3千語であった。出現語数に対する異なり語数の変化を図3に、出現頻度に対する累積出現率の変化を図4に示す。上位30位までの高頻出単語は表2に示す。

先に述べた衛星放送の分析と、AP電の調査から以下の特徴が判明している。

<衛星放送/AP電に共通>

- ・使われている用語が広範囲に及び、しかも分類が難しい。
- ・人名、機関名、国名、地名などの固有名詞が頻出する。人名は肩書を含んだ複雑な表現が多い。
- ・複雑な数量表現が頻出する。

address	内容	
(16進)		
00000	1c20 172f9	記事#1へのポインター
00010	1ca1 1849e	記事#2へのポインター
01c10	17225 d719f	最後の記事へのポインター
01c20	W1001---u IBX W0002 25-04 00756	記事#1のヘッダー
	233 87 10 lds ldp snp txm mon	
	"AP News Digest,0789<	
	"EDITORS:<	
	YY	記事#1のトレイラ
	"END<	
	AP-TK-25-04-90 0004GMT<	
	##	
01ca1	W1002---u SBX W0003 25-04 00384	記事#2のヘッダー
	233 87	
	"Nicaragua-Contras,0429<	
	"LaserPhotos available<	
	##	
172f9	MOSCOW - The KGB has reinforced its units on the Lithuanian border and increased patrols offshore, an off...	記事#1の本文
	By Mary McVean.	
	##	
1849e	Joaquina Garcia is a peasant woman who hasn't seen four sons since they left to join the Contra rebels nine	記事#2の本文
	##	
d7b33	##	

図2 英語ニュースデータベースのデータ形式

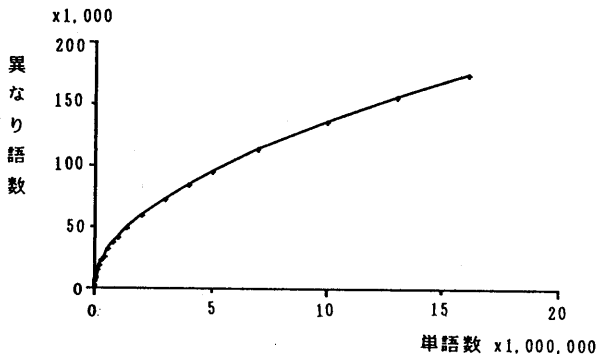


図3 出現語数に対する異なり語数の変化

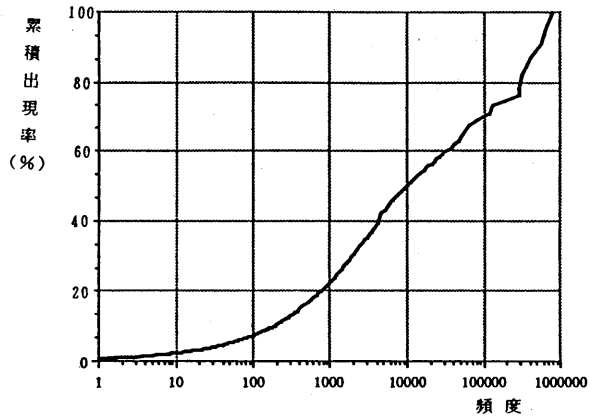


図4 出現頻度に対する語の累積出現率の変化

<AP電のニュース>

- ・1文が比較的長い。(23.7語/文)
- ・言動・思考に関する動詞が多用される。特に say は be動詞に次いで使われ have動詞より頻度が多い。call, talk, report, ask, think, want, feel の使用も多い。
- ・時制は「過去」の文が多い。

表2 AP電の高頻度単語(上位30位)  
(1600万語中, 単位千語)

the	982	on	109
, (カンマ)	749	he	89
, (ピリオド)	592	by	87
of	424	with	85
to	373	will	78
a	373	it	71
in	347	from	69
be動詞	344	at	68
and	328	as	65
数字	304	but	52
say	201	not	48
have動詞	159	his	47
's	137	government	44
for	132	who	42
that	123	U.S.	36

<衛星放送のニュース>

- ・ 1文が比較的短い。(11.0語/文)
- ・ 無生物主語の文が多い
- ・ 口語的な表現がよく使われる。
- ・ 時制は「現在」の文が多い。

書の構築」電子情報通信学会言語理解とコミュニケーション研究会, NLC91-6(1990)

- 6) 浦谷「高速な複数文字列照合アルゴリズム: F A S T」情処論文誌, Vol. 30, No. 9(1989)

5. おわりに

NHKの衛星放送で試用されている機械翻訳システムの現況について報告した。また、実際に機械翻訳に供された英語ニュース文の分析結果も示した。衛星放送のニュースで使用することによって生ずる点を中心にシステムの問題点について触れた。さらに、現在構築中の英語ニュースデータベースおよびシステム改善の方向性について述べた。

英語ニュースデータベースはやっと利用可能となったばかりであり、本格的な利用はこれからである。これを用いて用語、慣用句、文型の抽出を行い、辞書の充実と文法の改良を進めていく予定である。

翻訳ソフトウェアはカテナ・リソース研究所のSTARをベースに構築した。データ分析に当たって、同研究所の松田主任研究員に多大な協力を頂いたことを感謝します。

おわりに、研究の推進ならびに実用化に関して、ご指導を頂いた放送技術研究所の杉本所長、石田前次長、画像研究部の二宮部長、藤原前副部長に感謝します。

<参考文献>

- 1) 相沢ほか「衛星放送ワールドニュースの英日機械翻訳」情報処理学会第40回全国大会2F-1(1990)
- 2) 中瀬「英日機械翻訳システムにおける解析手法について」情報処理学会自然言語処理研究会, NL67-7(1988)
- 3) 加藤ほか「英日機械翻訳における固有名詞処理」情報処理学会第40回全国大会2F-2(1990)
- 4) 田中ほか「類語国語辞典を介した意味マーカー付与」情報処理学会第40回全国大会6F-7(1990)
- 5) 相沢ほか「英日機械翻訳のためのニュース用語辞